

MENU **SEARCH** **INDEX** **DETAIL** **JAPANESE** **BACK**

3 / 4

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 11-045268

(43)Date of publication of application : 16.02.1999

(51)Int.Cl. G06F 17/30
G06F 17/22
G06F 17/27
G06F 17/21

(21)Application number : 09-201985

(71)Applicant : JUST SYST CORP

(22)Date of filing : 28.07.1997

(72)Inventor : FUJITA SUMIO

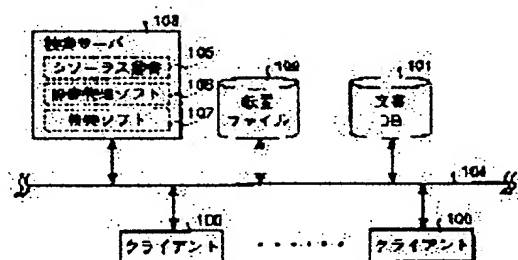
(54) DOCUMENT RETRIEVAL DEVICE AND COMPUTER-READABLE RECORDING MEDIUM WHERE EPROGRAM MAKING COMPUTER FUNTION AS SAME DEVICE IS RECORDED

(57)Abstract:

PROBLEM TO BE SOLVED: To reduce the trouble of thesaurus dictionary management by registering noun phrases as relative words and similar words of corresponding index words in a relative word and similar word dictionary.

SOLUTION: Clients 100 select index words in the thesaurus dictionary 105 and specify the execution of a retrieving process for documents. Dictionary management software 106 manages the thesaurus dictionary 105. Retrieval software 107 generates a dislocation file 102 by using documents in a document data base 101, uses a retrieval engine to retrieve a corresponding document from the dislocation file 102 according to a selected index word in the thesaurus dictionary 105, and extracts an index word or similar word of the thesaurus dictionary 105 from the document as the retrieval result. Namely, a score corresponding statistical information on the appearance frequency,

distribution, etc., of an object document group to be retrieved which is selected for the extracted noun phrase is imparted, the noun phrase having a score corresponding to a set retrieval condition is selected, and the noun phrase is registered as the relative word or synonym of the corresponding index word in the relative word or synonym.

**LEGAL STATUS**

[Date of request for examination]

[Date of sending the examiner's decision of rejection]

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平11-45268

(43) 公開日 平成11年(1999) 2月16日

(51) Int.Cl.⁶

識別記号

F I

G 0 6 F 17/30
17/22
17/27
17/21

G 0 6 F 15/40 3 7 0 A
15/20 5 2 2 L
5 5 0 F
5 7 0 N
15/40 3 7 0 J

審査請求 未請求 請求項の数 4 O L (全 11 頁) 最終頁に続く

(21) 出願番号 特願平9-201985

(22) 出願日 平成9年(1997) 7月28日

(71) 出願人 390024350

株式会社ジャストシステム
徳島県徳島市沖浜東3-46

(72) 発明者 藤田 澄男

徳島市沖浜東3丁目46番地 株式会社ジャ
ストシステム内

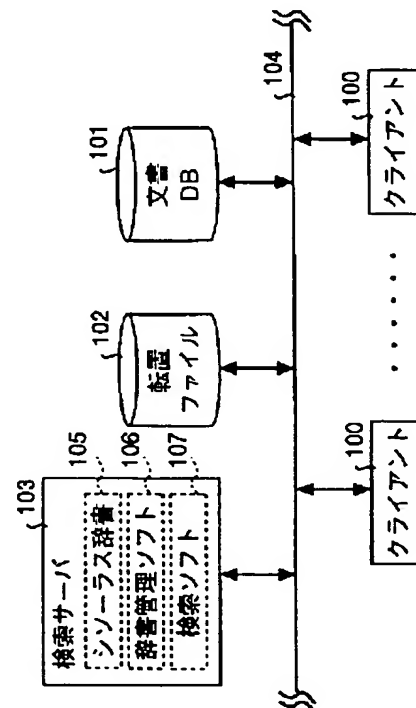
(74) 代理人 弁理士 酒井 昭徳

(54) 【発明の名称】 文書検索装置およびその装置としてコンピュータを機能させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体

(57) 【要約】

【課題】 検索によって得た文書から検索用のシソーラス辞書の索引語の関連語または類義語を自動的に抽出して登録できるようにすること。

【解決手段】 検索サーバ103は、意味体系に従って分類項目となる索引語を分類すると共に、関連語や類義語を該当する索引語にそれぞれ関連づけて記憶したシソーラス辞書105と、クライアント100で索引語が選択されると、選択された索引語およびその関連語や類義語を検索条件として検索処理を行うと共に、検索結果の文書から名詞句を抽出し、抽出した名詞句に対し、選択された文書および検索対象の文書群における出現頻度および分布等の統計情報に応じたスコアを付与し、予め設定された選択条件に該当するスコアの名詞句を選択する検索ソフト107と、選択した名詞句を該当する索引語の関連語または類義語としてシソーラス辞書105に登録する辞書管理ソフト106と、を備えている。



【特許請求の範囲】

【請求項 1】 検索条件に基づいて、検索対象の文書群から該当する文書を検索する文書検索装置において、意味体系に従って分類項目となる索引語を分類すると共に、前記索引語の関連語および／または類義語を前記分類した索引語にそれぞれ関連づけして記憶した関連語・類義語辞書と、前記関連語・類義語辞書の少なくとも索引語を画面表示する索引語表示手段と、前記索引語表示手段で画面表示された索引語を選択するための索引語選択手段と、前記索引語選択手段を介して索引語が選択されると、選択された索引語および前記索引語に関連づけられた関連語および／または類義語を前記検索条件として、該当する文書を検索する検索手段と、前記検索手段で検索した文書の一覧を表示する一覧表示手段と、前記一覧表示手段で表示された文書を選択するための文書選択手段と、前記文書選択手段を介して選択された文書から名詞句を抽出する名詞句抽出手段と、前記名詞句抽出手段で抽出した名詞句に対し、前記選択手段で選択した文書および検索対象の文書群における出現頻度および分布等の統計情報に応じたスコアを付与し、予め設定された選択条件に該当するスコアの名詞句を選択する名詞句選択手段と、前記名詞句選択手段で選択した名詞句を該当する索引語の関連語または類義語として前記関連語・類義語辞書に登録する辞書登録手段と、を備えたことを特徴とする文書検索装置。

【請求項 2】 前記辞書登録手段は、登録しようとする名詞句が既に該当する索引語の関連語または類義語として前記関連語・類義語辞書中に存在する場合、前記存在する関連語または類義語に正の重みを付与し、前記検索手段は、前記関連語または類義語に付与された重みを用いて、該当する文書の検索を行うことを特徴とする請求項 1 に記載の文書検索装置。

【請求項 3】 前記文書選択手段は、前記一覧表示手段で一覧表示された文書を選択する際に、前記索引語に適合する文書に対して正の重み付けを指定を行うことが可能であると共に、前記索引語に適合しない文書に対して負の重み付けの指定を行うことが可能であり、前記辞書登録手段は、前記正の重み付けが指定された文書から得た名詞句に正の重みを付与すると共に、前記負の重み付けが指定された文書から得た名詞句に負の重みを付与して該当する索引語の関連語または類義語として前記関連語・類義語辞書に登録し、前記検索手段は、前記関連語または類義語に付与された重みを用いて、該当する文書の検索を行うことを特徴とする請求項 1 または 2 に記載の文書検索装置。

【請求項 4】 前記請求項 1 ～ 3 のいずれか 1 つに記載の文書検索装置の各手段としてコンピュータを機能させるためのプログラムを記録したことを特徴とするコンピュータ読み取り可能な記録媒体。

【発明の詳細な説明】**【0001】**

【発明の属する技術分野】 本発明は、検索によって得た文書から検索用のシソーラス辞書の索引語の関連語または類義語を自動的に抽出して登録できるようにした文書検索装置およびその装置としてコンピュータを機能させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体に関する。

【0002】

【従来の技術】 複数の文書を格納した文書 DB（データベース）から特定の文書を検索する文書検索装置は、一般に、検索式や検索文等の検索条件を入力し、入力した検索条件に該当する文書を文書 DB から検索するものである。

【0003】 ところで、上記文書検索装置では、入力した検索条件に基づいて検索を行うため、検索条件の語彙そのものではなく、検索条件中の語彙に関連する語彙を用いて記述された文書については、入力した検索条件に該当せず、検索結果に漏れが生じることがあった。

【0004】 そこで、検索用のシソーラス辞書を予め用意しておき、入力した検索条件を索引語として、該当する関連語や類義語をシソーラス辞書から抽出し、入力した検索条件にシソーラス辞書から抽出した関連語や類義語を加えて検索を行うことにより、検索結果に漏れが生じることを防止した文書検索装置が提案されている。

【0005】

【発明が解決しようとする課題】 しかしながら、上記従来の技術においては、検索用のシソーラス辞書を人手で生成しているため、常に最新の類義語が登録されている状態でシソーラス辞書を管理することは困難であるという問題点があった。特に、先端技術に関しては、常に新しい用語が次々と生まれてくるため、新たな用語を収集し、シソーラス辞書に登録する作業を継続的に行うことは困難であった。そして、シソーラス辞書への登録を怠れば、いくらシソーラス辞書を用いて検索を行ったとしても、常に高い精度の検索結果を得ることは不可能であるという問題点があった。

【0006】 本発明は上記に鑑みてなされたものであって、検索によって得た文書から検索用のシソーラス辞書の索引語の関連語または類義語を自動的に抽出して登録できるようにすることにより、シソーラス辞書を管理するための労力の軽減を図ることを目的とする。

【0007】 また、本発明は上記に鑑みてなされたものであって、シソーラス辞書を常に最新の関連語および類義語が登録された状態に保つことができるようにすることにより、精度の高い検索結果を得ることができるよう

にすることを目的とする。

【0008】

【課題を解決するための手段】上記目的を達成するため、請求項1の文書検索装置は、検索条件に基づいて、検索対象の文書群から該当する文書を検索する文書検索装置において、意味体系に従って分類項目となる索引語を分類すると共に、前記索引語の関連語および／または類義語を前記分類した索引語にそれぞれ関連づけて記憶した関連語・類義語辞書と、前記関連語・類義語辞書の少なくとも索引語を画面表示する索引語表示手段と、前記索引語表示手段で画面表示された索引語を選択するための索引語選択手段と、前記索引語選択手段を介して索引語が選択されると、選択された索引語および前記索引語に関連づけられた関連語および／または類義語を前記検索条件として、該当する文書を検索する検索手段と、前記検索手段で検索した文書の一覧を表示する一覧表示手段と、前記一覧表示手段で表示された文書を選択するための文書選択手段と、前記文書選択手段を介して選択された文書から名詞句を抽出する名詞句抽出手段と、前記名詞句抽出手段で抽出した名詞句に対し、前記選択手段で選択した文書および検索対象の文書群における出現頻度および分布等の統計情報に応じたスコアを付与し、予め設定された選択条件に該当するスコアの名詞句を選択する名詞句選択手段と、前記名詞句選択手段で選択した名詞句を該当する索引語の関連語または類義語として前記関連語・類義語辞書に登録する辞書登録手段と、を備えたものである。

【0009】また、請求項2の文書検索装置は、請求項1に記載の文書検索装置において、前記辞書登録手段が、登録しようとする名詞句が既に該当する索引語の関連語または類義語として前記関連語・類義語辞書中に存在する場合、前記存在する関連語または類義語に正の重みを付与し、前記検索手段が、前記関連語または類義語に付与された重みを用いて、該当する文書の検索を行うものである。

【0010】また、請求項3の文書検索装置は、請求項1または2に記載の文書検索装置において、前記文書選択手段が、前記一覧表示手段で一覧表示された文書を選択する際に、前記索引語に適合する文書に対して正の重み付けを指定を行うことが可能であると共に、前記索引語に適合しない文書に対して負の重み付けの指定を行うことが可能であり、前記辞書登録手段が、前記正の重み付けが指定された文書から得た名詞句に正の重みを付与すると共に、前記負の重み付けが指定された文書から得た名詞句に負の重みを付与して該当する索引語の関連語または類義語として前記関連語・類義語辞書に登録し、前記検索手段が、前記関連語または類義語に付与された重みを用いて、該当する文書の検索を行うものである。

【0011】さらに、請求項4のコンピュータ読み取り可能な記録媒体は、前記請求項1～3のいずれか1つに

記載の文書検索装置の各手段としてコンピュータを機能させるためのプログラムを記録したものである。

【0012】

【発明の実施の形態】以下、本発明の文書検索装置およびその装置としてコンピュータを機能させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体の一実施の形態について、添付の図面を参照しつつ詳細に説明する。

【0013】図1は、本実施の形態の文書検索装置のシステム構成図である。図1に示す文書検索装置は、後に詳細に説明するシソーラス辞書105中の索引語を選択することにより、文書の検索処理の実行を指定する機能を有する複数のクライアント100と、シソーラス辞書105、シソーラス辞書105を管理する辞書管理ソフト106、および文書DB（データベース）101中の文書を用いて転置ファイル102を生成すると共に、ベクトル空間法を利用した検索エンジン（例えば、CLARITECH社のCLARIT等）を用いることにより、クライアント100で選択されたシソーラス辞書105中の索引語に基づいて、転置ファイル102から該当する文書を検索し、さらに、検索結果の文書からシソーラス辞書105の索引語の関連語または類義語を抽出する検索ソフト107を備えた検索サーバ103と、上記クライアント100や検索サーバ103等を接続するネットワーク104と、から構成されている。

【0014】図1において、文書DB101は、クライアント100等で作成された複数の文書を格納したものであり、格納される文書は、ワープロ文書や、SGML、HTML等の構造化文書等、いかなる種類の文書であっても良い。本実施の形態においては、文書DB101に格納された文書を検索対象とするが、検索対象を文書DB101中の文書に限定するものではない。

【0015】転置ファイル102は、文書DB101中の複数の文書と、これら複数の文書それぞれから後述する方法で抽出した複数の索引語との関係を規定することにより、ある索引語が各文書それぞれにおいてどの程度重要な語彙であるかをベクター表現を用いて示したものであって、この索引語を用いて該当する文書を検索することができるようにしたものである。

【0016】具体的には、1つの文書を複数のセンテンスからなるサブドキュメント単位に区切り、サブドキュメントから上記索引語となる名詞句を抽出して、抽出した名詞句それぞれについて、サブドキュメント中の出現頻度、文書DB101全体における分布等の統計情報を求め、求めた名詞句毎の統計情報を用いて各サブドキュメントをベクター表現に変換する。そして、変換したサブドキュメントのベクター表現に基づいて、文書のベクター表現を生成する。転置ファイル102は、このようにしてベクター表現された文書DB101中の文書を格納するものである。

【0017】なお、各索引語には、対応する文書中の重要度に応じた重み付けを行うことができる。また、文書のベクター表現については、実際の検索を行う際に、サブドキュメントのベクター表現に基づいて生成することにしても良い。

【0018】クライアント100および検索サーバ103は、パーソナルコンピュータやワークステーション等によって構成される。

【0019】図2は、検索サーバ103におけるシソーラス辞書105の内容を画面表示した様子の一例を示す説明図である。シソーラス辞書105は、意味体系に従って分類項目となる索引語を分類すると共に、索引語の関連語および／または類義語を分類した索引語にそれぞれ関連づけして記憶したものである。図2に示したシソーラス辞書105は、例として、新聞記事の記事分類に従い、階層構造となるように索引語を分類したものであり、フォルダ（またはディレクトリ）名のようにして表示されているものが索引語に該当する。また、ある索引語とその下位にある索引語とは、下位の索引語が上位の索引語の関連語または類義語となる。

【0020】図2において、索引語「交通事故」について考えると、索引語「交通事故」は、索引語「社会面記事」の関連語となり（図2の左側の部分を参照）、下位の索引語「人身事故」、「物損事故」、「保険」は、それぞれ索引語「交通事故」の関連語または類義語となっている（図2の右側の部分を参照）。さらに、索引語「交通事故」には、図2の右側の部分に示すように、「衝突」、「死傷者」、「追突」、「脇見」、「飲酒」、「業務上過失致死」等が関連語または類義語として関連づけられている。なお、関連語および類義語には、索引語に対する関連性や類似性に応じて、それぞれ重みが付されており、付された重みを文書の検索の際に用いることができるようになっている。

【0021】このシソーラス辞書105は、辞書管理ソフト106を介してクライアント100からアクセスできるようになっており、図2は、クライアント100でシソーラス辞書105の内容を画面表示した様子である。クライアント100においては、検索したい文書が該当する索引語を探し、その索引語を図示しないマウス等で選択することにより、選択した索引語に関連づけられた関連語および／または類義語を用いた検索の実行を指定することができる。

【0022】また、図3は、検索サーバ103において、検索ソフト107の処理を示す概略ブロック図である。検索ソフト107は、文書DB101中の文書を転置ファイル102に登録する処理と、ベクトル空間法を利用した検索処理と、検索結果の文書からシソーラス辞書105中の索引語の関連語または類義語を抽出する処理を行うものである。

【0023】検索ソフト107において、転置ファイル

102に登録する処理は、自然言語処理モジュール300と、データベース・ビルド・コンポーネント304とによって行われる。

【0024】具体的に、自然言語処理モジュール300は、文書DB101から文書を入力し、文書のフォーマットの認識処理や、品詞情報等を格納した辞書301および各単語の係り受け等を解析するための文法辞書302を用いて形態素解析、構文解析、名詞句抽出等の解析処理を行い、上述したサブドキュメント毎の名詞句リストを含むドキュメント・セット303を生成する。

【0025】データベース・ビルドコンポーネント304は、自然言語処理モジュール300で生成したドキュメント・セット303を入力し、入力したドキュメント・セット303中の各サブドキュメントを上述したようにしてベクター表現に変換すると共に、サブドキュメントのベクター表現に基づいて、文書のベクター表現を生成して転置ファイル102に登録する。

【0026】また、検索ソフト107において、文書の検索処理は、自然言語処理モジュール300と、クエリー・ビルド・コンポーネント305と、検索エンジン307とによって行われる。

【0027】具体的に、自然言語処理モジュール300は、辞書管理ソフト106を介して、クライアント100で選択されたシソーラス辞書105中の索引語とその索引語に関連づけられた関連語および／または類義語を検索条件として入力し、入力した検索条件の索引語・関連語または類義語毎にドキュメント・セット303を生成する。

【0028】クエリー・ビルド・コンポーネント305は、ドキュメント・セット303を入力し、索引語、関連語または類義語について、文書DB101（転置ファイル102）全体における分布等の統計情報を求め、求めた統計情報と上記重みを用いてドキュメント・セット303をベクター表現に変換したクエリー・ドキュメント306を生成する。

【0029】検索エンジン307は、クエリー・ビルド・コンポーネント305で生成したクエリー・ドキュメント306を入力し、転置ファイル102中の各文書のベクター表現とクエリー・ドキュメント306（索引語、関連語または類義語のベクター表現）とを比較して、クエリー・ドキュメント306との類似度に応じたスコアを各文書に付与し、所定の閾値を超えるスコアが付与された文書リスト308を検索結果として出力する。

【0030】さらに、検索ソフト107において、シソーラス辞書105に登録する関連語または類義語の抽出処理は、自然言語処理モジュール300と、シソーラス抽出エンジン309とによって行われる。

【0031】具体的に、自然言語処理モジュール300は、上述した検索処理の結果に基づいて、クライアント

100で選択された文書を文書DB101から入力し、入力した文書について、フォーマットの認識処理や、品詞情報等を格納した辞書301および各単語の係り受け等を解析するための文法辞書302を用いて形態素解析、構文解析、名詞句抽出等の解析処理を行い、上述したサブドキュメント毎の名詞句リストを含むドキュメント・セット303を生成する。

【0032】シソーラス抽出エンジン310は、自然言語処理モジュール300で生成したドキュメント・セット303を入力し、入力したドキュメント・セット303中の各名詞句それぞれについて、各文書（ドキュメント・セット303）中の出現頻度や文書DB101（転置ファイル102）中の分布等の統計データを演算し、演算した統計データに基づいて、各名詞句にスコアを付与する。そして、予め設定した閾値を超えるスコアの名詞句を選択してシソーラスリスト310を生成し、辞書管理ソフト106に出力する。

【0033】そして、辞書管理ソフト106は、検索ソフト103からシソーラスリスト310を入力し、シソーラスリスト310中の名詞句を、検索を開始する際にクライアント100で選択された索引語の関連語または類義語としてシソーラス辞書105に登録する。

【0034】なお、図1においては、文書DB101および転置ファイル102をネットワーク104に単独に接続した構成を示したが、これらを検索サーバ103に直接接続する構成としても良い。また、図1においては、本実施の形態の文書検索装置をネットワーク104を介したシステムで構成するように示したが、クライアント100と検索サーバ103の処理を1つのコンピュータで行うようにすることもできる。

【0035】次に、上述した構成を備えた文書検索装置の動作について、（1）転置ファイルの生成処理、

（2）文書の検索処理、（3）シソーラス辞書への登録処理の順で詳細に説明する。

【0036】（1）転置ファイルの生成処理

図4は、転置ファイルの生成処理を示すフローチャートである。検索サーバ103は、新たな文書が文書DB101に登録された場合（S401）、この文書を入力して転置ファイル102に登録するための処理を開始する（S402）。

【0037】検索サーバ103において、自然言語処理モジュール300は、ステップS402で入力した文書を解析する処理を行う（S403）。具体的には、入力した文書がワープロ文書、HTML等の構造化文書等、いかなるフォーマットの文書であるかを判定する処理を行う。その後、辞書301および文法辞書302を用いて形態素解析、係り受け等の構文解析を行い、文書を複数のサブドキュメントに区分すると共に、区分したサブドキュメントから名詞句を抽出する等の処理を行う。

【0038】そして、自然言語処理モジュール300

は、ステップS403における処理の結果に基づいて、サブドキュメント毎に名詞句リストを生成し、生成した名詞句リストを含むドキュメント・セット303を生成する（S404）。

【0039】その後、データベース・ビルド・コンポーネント304は、自然言語処理モジュール300で生成したドキュメント・セット303を入力し、文書のベクター表現を生成して転置ファイル102に登録する処理を行う（S405）。

【0040】具体的には、ドキュメント・セット303中のサブドキュメントの各名詞句を転置ファイル102の索引語として、サブドキュメント中の出現頻度、文書DB101全体における分布等の統計情報を求め、求めた名詞句毎の統計情報を用いてサブドキュメントをベクター表現に変換する。この処理をドキュメント・セット303中の全てのサブドキュメントについて行い、変換したサブドキュメントのベクター表現に基づいて、文書のベクター表現を生成して転置ファイル102に登録する。その結果、文書DB101に新たに登録された文書が転置ファイル102に登録されることになる。

【0041】（2）文書の検索処理

続いて、上述したようにして生成した転置ファイル102に基づいて、文書DB101中から特定の文書を検索するための処理について説明する。図5は、文書の検索処理を示すフローチャートである。

【0042】ユーザは、クライアント100を操作して、検索サーバ103の辞書管理ソフト107にシソーラス辞書105の内容の表示を要求する。その結果、クライアント100に図2に示したシソーラス辞書105の内容が画面表示される。

【0043】そこで、ユーザは、シソーラス辞書105中の索引語を参照し、検索によって得たい文書が該当する索引語を選択する。すなわち、この索引語は、検索条件の役割を果たすことになる。辞書管理ソフト107は、クライアント100で索引語が選択されると（S501）、該当する関連語および／または類義語を検索条件として検索ソフト107に出力する（S502）。

【0044】例えば、図2に示した「交通事故」という索引語がクライアント100によって選択されたとすると、辞書管理ソフト106は、索引語「交通事故」と、索引語「交通事後」の関連語および／または類義語を検索条件として検索ソフト107に出力することになる。

【0045】検索ソフト107は、辞書管理ソフト106から検索条件を入力し、自然言語処理モジュール300において、検索条件を構成する上記索引語、関連語および／または類義語毎にドキュメント・セット303を生成する（S503）。

【0046】続いて、クエリー・ビルド・コンポーネント305は、自然言語処理モジュール300からドキュメント・セット303を入力し、索引語、関連語または

類義語について、文書DB101（転置ファイル102）全体における分布等の統計情報を求め、求めた統計情報とそれらに付与された重みを用いてドキュメント・セット303をベクター表現に変換したクエリー・ドキュメント306を生成する（S504）。

【0047】検索エンジン307は、クエリー・ビルド・コンポーネント305で生成したクエリー・ドキュメント306を入力し、転置ファイル102中の各文書のベクター表現とクエリー・ドキュメント306（索引語、関連語または類義語のベクター表現）とを比較して、クエリー・ドキュメント306との類似度に応じたスコアを各文書に付与する（S505）。すなわち、ベクトル空間法を用いた検索処理が行われる。

【0048】なお、類似度に応じたスコアは、各文書とクエリー・ドキュメント306との類似度を余弦距離に基づいて表現したものであり、スコアが大きい文書がよりクエリー・ドキュメント306と類似していることを表している。

【0049】そして、検索エンジン307は、予め設定されたスコアの閾値に基づいて、閾値を超えるスコアが付与された文書を選択し、選択した文書に基づいて、文書リスト308を生成し、クライアント100に出力する（S506）。

【0050】クライアント100は、検索サーバ103から文書リスト308を入力し、入力した文書リスト308に基づいて、上位のランキングの文書から順に、該当する索引語に関連づけて文書一覧を表示する（S507）。クライアント100のユーザは、一覧表示された文書から所望の文書を選択することにより、その文書を画面表示することができる。

【0051】（3）シソーラス辞書への登録処理
さらに、検索結果の文書から名詞句を抽出してシソーラス辞書へ登録する処理について説明する。図6は、シソーラス辞書への登録処理を示すフローチャートである。

【0052】クライアント100のユーザは、画面表示された文書一覧から検索結果としてふさわしい文書（選択したシソーラス辞書105中の索引語に適合する文書）を選択し、選択した文書を検索結果として検索サーバ103に出力する。検索サーバ103の検索ソフト105は、クライアント100から検索結果を入力すると（S601）、検索結果に該当する文書を文書DB101から入力する（S602）。

【0053】文書DB101から文書を入力すると、自然言語処理モジュール200は、入力した文書毎に、フォーマットの認識処理や、品詞情報等を格納した辞書201および各単語の係り受け等を解析するための文法辞書202を用いて形態素解析、構文解析、名詞句抽出等の解析処理を行う（S603）。

【0054】その後、ステップS603における解析処理の結果に基づいて、サブドキュメント毎の名詞句リス

トを含むドキュメント・セット204を1文書を単位として生成する（S604）。

【0055】シソーラス抽出エンジン309は、自然言語処理モジュール200で生成したドキュメント・セット204を入力し、入力したドキュメント・セット204中の各名詞句それぞれについて、各文書（ドキュメント・セット204）中の出現頻度や文書DB101（転置ファイル102）中の分布等の統計データを演算する（S605）。

【0056】ステップS605で統計データを演算した後、シソーラス抽出エンジン309は、求めた統計データに基づいて、各名詞句に対してスコア付けを行う（S606）。このスコアは、文書における各名詞句の重要性および検索を行う際に選択された索引語に対する関連性または類似性を表すもので、スコアが大きいもの程、重要性および関連性または類似性が高いことを表している。

【0057】シソーラス抽出エンジン309は、ステップS606で行ったスコア付けの結果に基づいて、予め設定された閾値を超えるスコアの名詞句を、クライアント100で選択された索引語（図5のステップS501参照）の関連語または類義語として抽出する（S607）。なお、ここでは、名詞句を抽出する条件として閾値を用いることにしたが、閾値に代えて、例えば、上位5番までのスコアの名詞句を抽出することにしても良い。

【0058】その後、シソーラス抽出エンジン309は、ステップS607で抽出した名詞句のリストであるシソーラスリスト310を生成して、辞書管理ソフト106に出力する（S608）。

【0059】辞書管理ソフト106は、検索ソフト107からシソーラスリスト310を入力し、入力したシソーラスリスト310中の名詞句を該当する索引語の関連語または類義語として、シソーラス辞書105に登録する（S609）。

【0060】なお、辞書管理ソフト106は、登録しようとする名詞句が該当する索引語の関連語または類義語として既にシソーラス辞書105中に存在する場合には、存在する関連語または類義語に正の重みを付与する。したがって、検索を行う際に、これらの関連語または類義語を含む文書がヒットする率が高められる。

【0061】また、検索された文書が一覧表示され、クライアント100で一覧表示された文書から索引語に適合する文書を選択する際には、索引語に適合する文書に対して正の重み付けを指定を行うことができると共に、索引語に適合しない文書に対して負の重み付けの指定を行うことができる。そして、辞書管理ソフト106は、シソーラスリスト310の名詞句をシソーラス辞書105に登録する際に、正の重み付けが指定された文書から得た名詞句には正の重みを付与して関連語または類義語

として登録すると共に、負の重み付けが指定された文書から得た名詞句には負の重みを付与して関連語または類義語として登録する。その結果、正の重みが付与された関連語または類義語を含む文書が検索でヒットする確率が高まり、一方、負の重みが付与された関連語または類義語を含む文書が検索でヒットする確率は低められることになる。なお、検索結果の文書から得た名詞句のうち、正の重みの指定がなされた文書と負の重みの指定がなされた文書の両方に存在する名詞句については、シソーラス辞書 105 に登録しないようにすることもできる。

【0062】さらに、図6のステップS609において、生成したシソーラスリスト310中の名詞句をそのままシソーラス辞書105に登録することにしたが、これらの名詞句を一度クライアント100に画面表示し、ユーザが選択した名詞句のみをシソーラス辞書105に登録することにしても良い。

【0063】このように、本実施の形態の文書検索装置によれば、検索によって得た文書からシソーラス辞書105中の索引語の関連語または類義語を自動的に抽出してシソーラス辞書105に登録できるようにすることにしたため、シソーラス辞書105を管理するための労力の軽減を図ることができると共に、シソーラス辞書105を常に最新の関連語や類義語が登録された状態に保つことができる。したがって、精度の高い検索処理を行うことができる。

【0064】なお、本実施の形態においては、ベクトル空間法による検索を例にとって説明したが、ブーリアン検索により検索処理を行うことにしても良い。

【0065】また、本実施の形態で説明した文書検索装置は、予め用意されたプログラムをコンピュータやワークステーションで実行することによって実現される。このプログラムは、ハードディスク、フロッピーディスク、CD-ROM、MO、DVD等のコンピュータで読み取り可能な記録媒体に記録され、コンピュータによって記録媒体から読み出されることによって実行される。また、このプログラムは、上記記録媒体を介して、またはネットワークを介して配布することができる。

【0066】

【発明の効果】以上説明したように、本発明の文書検索装置（請求項1）によれば、意味体系に従って分類項目となる索引語を分類すると共に、索引語の関連語および／または類義語を分類した索引語にそれぞれ関連づけして記憶した関連語・類義語辞書と、関連語・類義語辞書の少なくとも索引語を画面表示する索引語表示手段と、索引語表示手段で画面表示された索引語を選択するための索引語選択手段と、索引語選択手段を介して索引語が選択されると、選択された索引語および索引語に関連づけられた関連語および／または類義語を検索条件として、該当する文書を検索する検索手段と、検索手段で検

索した文書の一覧を表示する一覧表示手段と、一覧表示手段で表示された文書を選択するための文書選択手段と、文書選択手段を介して選択された文書から名詞句を抽出する名詞句抽出手段と、名詞句抽出手段で抽出した名詞句に対し、選択手段で選択した文書および検索対象の文書群における出現頻度および分布等の統計情報に応じたスコアを付与し、予め設定された選択条件に該当するスコアの名詞句を選択する名詞句選択手段と、名詞句選択手段で選択した名詞句を該当する索引語の関連語または類義語として関連語・類義語辞書に登録する辞書登録手段と、を備えたため、関連語・類義語辞書を管理するための労力の軽減を図ることができると共に、関連語・類義語辞書を常に最新の関連語や類義語が登録された状態に保つことができる。したがって、精度の高い検索処理を行うことができる。

【0067】また、本発明の文書検索装置（請求項2）によれば、請求項1に記載の文書検索装置において、辞書登録手段は、登録しようとする名詞句が既に該当する索引語の関連語または類義語として関連語・類義語辞書中に存在する場合、存在する関連語または類義語に正の重みを付与し、検索手段は、関連語または類義語に付与された重みを用いて、該当する文書の検索を行うため、検索結果に応じて関連語・類義語辞書を学習させることができ、検索を重ねる毎に検索精度の向上を図ることができる。

【0068】また、本発明の文書検索装置（請求項3）によれば、請求項1または2に記載の文書検索装置において、文書選択手段は、一覧表示手段で一覧表示された文書を選択する際に、索引語に適合する文書に対して正の重み付けを指定を行うことが可能であると共に、索引語に適合しない文書に対して負の重み付けの指定を行うことが可能であり、辞書登録手段は、正の重み付けが指定された文書から得た名詞句に正の重みを付与すると共に、負の重み付けが指定された文書から得た名詞句に負の重みを付与して該当する索引語の関連語または類義語として関連語・類義語辞書に登録し、検索手段は、関連語または類義語に付与された重みを用いて、該当する文書の検索を行うため、検索結果に応じて関連語・類義語辞書を学習させることができ、検索を重ねる毎に検索精度の向上を図ることができる。

【0069】さらに、本発明のコンピュータ読み取り可能な記録媒体（請求項4）によれば、請求項1～3のいずれか1つに記載の文書検索装置の各手段としてコンピュータを機能させるためのプログラムを記録したため、このプログラムをコンピュータに実行させることにより、関連語・類義語辞書を管理するための労力の軽減を図ることができると共に、関連語・類義語辞書を常に最新の関連語や類義語が登録された状態に保つことができ、精度の高い検索処理を行うことができる文書検索装置を提供することができる。

【図面の簡単な説明】

【図 1】 本実施の形態の文書検索装置のシステム構成図である。

【図 2】 本実施の形態の文書検索装置において、検索サーバにおけるシソーラス辞書の内容を画面表示した様子の一例を示す説明図である。

【図 3】 本実施の形態の文書検索装置において、検索サーバにおける検索ソフトの処理を示す概略ブロック図である。

【図 4】 本実施の形態の文書検索装置において、転置ファイルの生成処理を示すフローチャートである。

【図 5】 本実施の形態の文書検索装置において、文書の検索処理を示すフローチャートである。

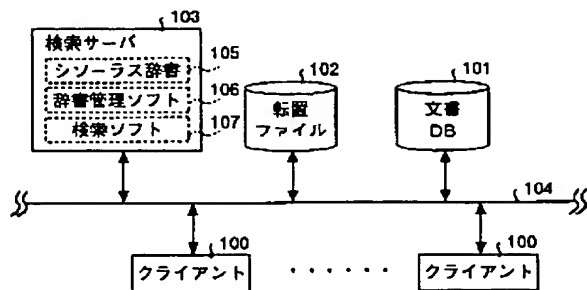
【図 6】 本実施の形態の文書検索装置において、シソーラス辞書への登録処理を示すフローチャートである。

【符号の説明】

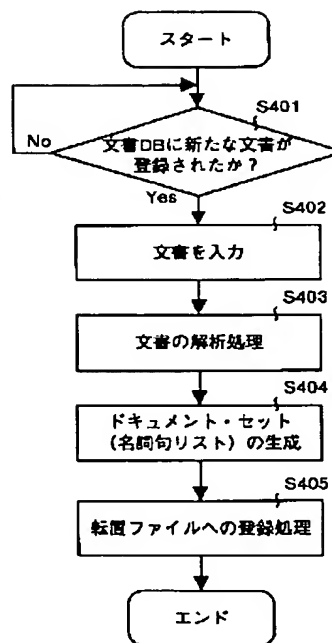
100 クライアント
101 文書DB

102 転置ファイル
103 検索サーバ
104 ネットワーク
105 シソーラス辞書
106 辞書管理ソフト
107 検索ソフト
300 自然言語処理モジュール
301 辞書
302 文法辞書
303 ドキュメント・セット
304 データベース・ビルド・コンポーネント
305 クエリー・ビルド・コンポーネント
306 クエリー・ドキュメント
307 検索エンジン
308 文書リスト
309 シソーラス抽出エンジン
310 シソーラスリスト

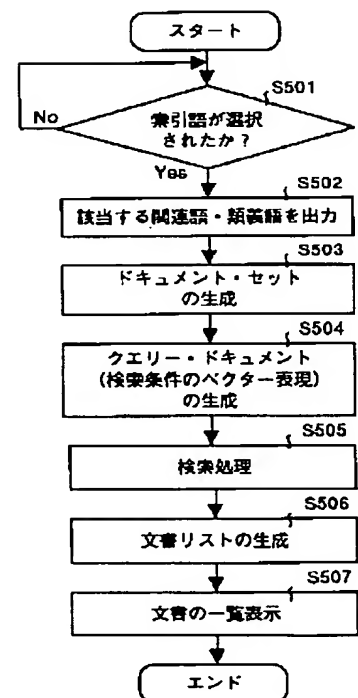
【図 1】



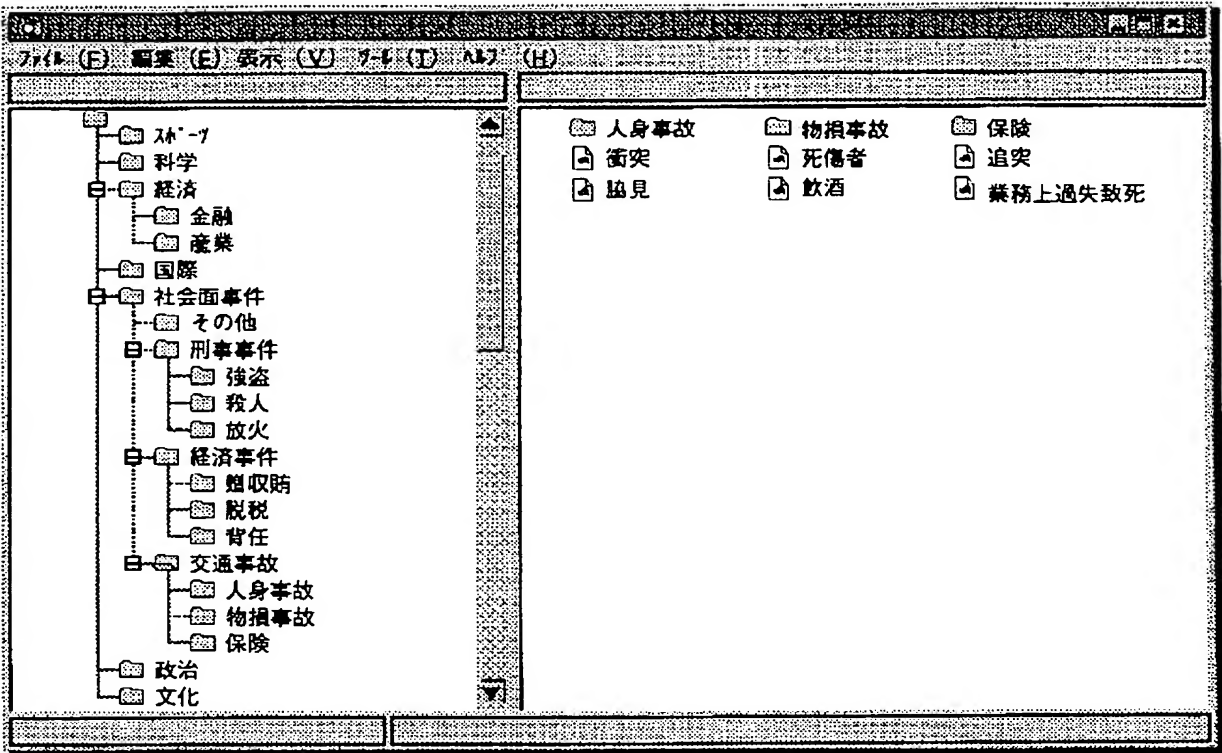
【図 4】



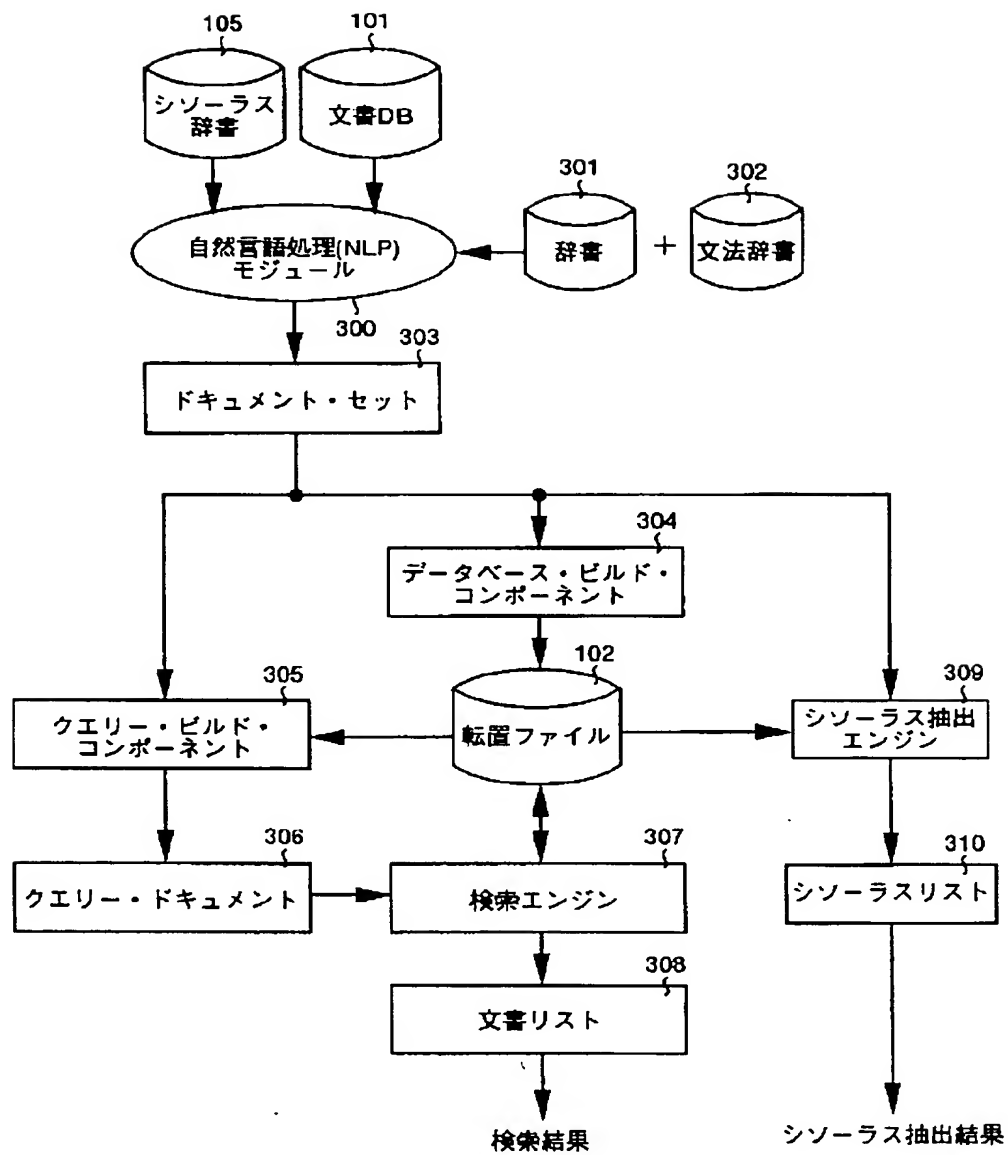
【図 5】



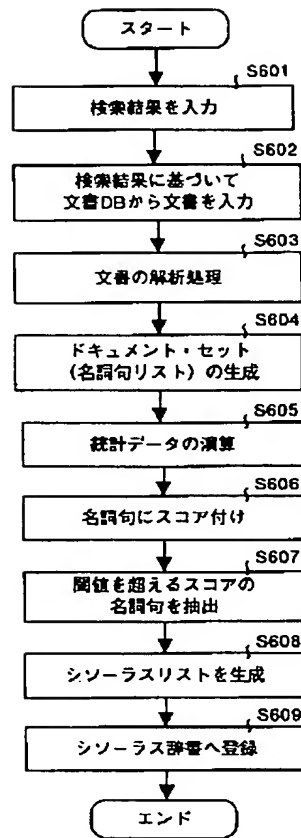
【図 2】



【図 3】



【図 6】



フロントページの続き

(51) Int. Cl. 6

識別記号

F I

G O 6 F 15/403

3 2 0 D